

A feature selection method using fixed-point algorithm for DNA microarray gene expression data

Alok Sharma^{a,b,c,*}, Kuldip K. Paliwal^b, Seiya Imoto^a, Satoru Miyano^a, Vandana Sharma^d and Rajeshkannan Ananthanarayanan^c

^aLaboratory of DNA Information Analysis, Human Genome Center, Institute of Medical Science, University of Tokyo, Tokyo, Japan

^bSchool of Engineering, Griffith University, Tokyo, Japan

^cSchool of Engineering and Physics, University of the South Pacific, Tokyo, Japan

^dFiji School of Medicine, University of the South Pacific, Tokyo, Japan

Abstract. As the performance of hardware is limited, the focus has been to develop objective, optimized and computationally efficient algorithms for a given task. To this extent, fixed-point and approximate algorithms have been developed and successfully applied in many areas of research. In this paper we propose a feature selection method based on fixed-point algorithm and show its application in the field of human cancer classification using DNA microarray gene expression data. In the fixed-point algorithm, we utilize between-class scatter matrix to compute the leading eigenvector. This eigenvector has been used to select genes. In the computation of the eigenvector, the eigenvalue decomposition of the scatter matrix is not required which significantly reduces its computational complexity and memory requirement.

Keywords: Feature selection, fixed-point algorithm, DNA microarray gene expression data, fast PCA

1. Introduction

Fixed-point algorithms have been recently applied to do many important applications such as independent component analysis (ICA) [9] and principal component analysis (PCA) [16]. Their popularity [1,3–5,8,11,13–15,19,22–24] is due to many reasons like low cost hardware implementation, low memory requirement, less processing time and less computational complexity.

In this paper, we propose a feature selection method using fixed-point algorithm of PCA [14]. The fixed-

point algorithm of PCA is also known as fast PCA (FPCA) algorithm. The FPCA algorithm has been recently extended and applied in face recognition [8,14,15], communication [13,24], VLSI architecture design [3–5] and in other areas or applications like in Yang et al. [23]; Shi and Guo [19]; Lai and Huang [11]; Wang et al. [22]; Albanese et al. [1]. The feature selection method plays a significant role in identifying crucial genes related to human cancers. It helps in understanding the gene regulation mechanism of cancer heterogeneity. We have carried out gene selection on DNA microarray gene expression datasets. These datasets, consisting of several thousands of gene expression profiles, have been widely used in the past for cancer classification problem. The fixed-point algorithm for feature selection has been proposed for lower computational time and memory requirement.

*Corresponding author: Alok Sharma, Laboratory of DNA Information Analysis, Human Genome Center, Institute of Medical Science, University of Tokyo, 4-6-1 Shirokanedai, Minato-ku, Tokyo, 108-8639, Japan. Tel.: +81 3 5449 5615; Fax: +81 3 5449 5442; E-mail: aloks@ims.u-tokyo.ac.jp.

FPCA algorithm works on training data $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ in a non-supervised manner. It does not require class labels for individual feature vectors \mathbf{x}_j . In FPCA algorithm, an eigenvector is computed by iteratively multiplying a covariance matrix $\Sigma_{\mathbf{x}} = \mathbf{H}\mathbf{H}^T$ (where $\mathbf{H} = \frac{1}{\sqrt{n}}[(\mathbf{x}_1 - \boldsymbol{\mu}), (\mathbf{x}_2 - \boldsymbol{\mu}), \dots, (\mathbf{x}_n - \boldsymbol{\mu})]$ and $\boldsymbol{\mu} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$ is the centroid of training data) with a random vector $\boldsymbol{\Phi} \in \mathbb{R}^{d \times 1}$ and updating $\boldsymbol{\Phi}$ as $\boldsymbol{\Phi} \leftarrow \Sigma_{\mathbf{x}} \boldsymbol{\Phi}$. The iteration process is terminated if some error criterion is below the threshold value.

The computational complexity of computing covariance matrix $\Sigma_{\mathbf{x}}$ is $O(d^2n)$. If the size of data dimensionality is very large then explicitly computing the covariance matrix $\Sigma_{\mathbf{x}}$ would be expensive. In that case, the updating can be done in the following manner: instead of computing $\Sigma_{\mathbf{x}}$ explicitly, a vector $\mathbf{g} = \mathbf{H}^T \boldsymbol{\Phi}$ can be computed first, and then $\boldsymbol{\Phi}$ can be updated as $\boldsymbol{\Phi} \leftarrow \mathbf{H}\mathbf{g}$. The computation of vector \mathbf{g} would require $2dn$ flops and the computation of $\boldsymbol{\Phi}$ using the product $\mathbf{H}\mathbf{g}$ would require $2dn$ flops. Therefore, the total flops to compute $\boldsymbol{\Phi}$ is $4dn$ per iteration.

It is known that both the range space and null space of between-class scatter matrix, \mathbf{S}_B , contain significant discriminant information [18]. Therefore, we use \mathbf{S}_B matrix by replacing $\Sigma_{\mathbf{x}}$ in the FPCA algorithm. We compute the leading eigenvector recursively until the desired number of genes is selected. We have compared the proposed method with other feature selection methods and promising results have been obtained. Since FPCA algorithm has been used in our strategy, we do not require to perform EVD of a matrix. This reduces the computational complexity and memory requirement significantly. Our method is, therefore, suited for a low cost hardware system.

2. Basic descriptions

In this section we describe the basic notations used in the paper. Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$ denote n training samples (or feature vectors) in a d -dimensional space having class labels $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$, where $\omega \in \{1, 2, \dots, c\}$ and c are the number of classes. The dataset \mathbf{X} can be subdivided into c subsets $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_c$, where \mathbf{X}_j belongs to class j and consists of n_j number of samples such that $n = \sum_{j=1}^c n_j$. The data subset $\mathbf{X}_j \subset \mathbf{X}$ and $\mathbf{X}_1 \cup \mathbf{X}_2 \cup \dots \cup \mathbf{X}_c = \mathbf{X}$. If $\boldsymbol{\mu}_j$ is the centroid of \mathbf{X}_j and $\boldsymbol{\mu}$ is the centroid of \mathbf{X} , then the between-class scatter matrix \mathbf{S}_B is defined as [6, 17]

$$\mathbf{S}_B = \sum_{j=1}^c n_j (\boldsymbol{\mu}_j - \boldsymbol{\mu})(\boldsymbol{\mu}_j - \boldsymbol{\mu})^T,$$

where

$$\boldsymbol{\mu}_j = \frac{1}{n_j} \sum_{\mathbf{x} \in \mathbf{X}_j} \mathbf{x}$$

and

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{\mathbf{x} \in \mathbf{X}} \mathbf{x}.$$

The between-class scatter matrix is a positive-semidefinite symmetric matrix which can be formed by using rectangular matrix; i.e., $\mathbf{S}_B = \mathbf{B}\mathbf{B}^T$, where rectangular matrix $\mathbf{B} \in \mathbb{R}^{d \times c}$ can be defined as [18]

$$\mathbf{B} = [\sqrt{n_1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}), \sqrt{n_2}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}), \dots, \sqrt{n_c}(\boldsymbol{\mu}_c - \boldsymbol{\mu})]. \quad (1)$$

3. The fixed-point algorithm for gene selection

The between-class scatter matrix \mathbf{S}_B contains significant discriminant information for classification [12]. We utilize \mathbf{S}_B and apply it in the framework of FPCA. The obtained orientation matrix \mathbf{W} from this procedure will be orthogonal. However, we are interested only in the leading eigenvector for gene selection. We modify the step of $\boldsymbol{\Phi} \leftarrow \Sigma_{\mathbf{x}} \boldsymbol{\Phi}$ of FPCA procedure as follows:

$$\mathbf{w} \leftarrow \mathbf{S}_B \mathbf{w}, \quad (2)$$

$$\mathbf{w} \leftarrow \text{orthonormalize}(\mathbf{w}). \quad (3)$$

Since for DNA microarray gene expression data, the size of \mathbf{S}_B matrix will be too large (as $d \gg n$), we can update Eq. (2) into two steps as:

$$\mathbf{w} \leftarrow \mathbf{B}^T \mathbf{w}, \quad (4)$$

$$\mathbf{w} \leftarrow \mathbf{B}\mathbf{w}. \quad (5)$$

If we define the fixed-point algorithm of Sharma and Paliwal [16] as $\boldsymbol{\Phi}_j \leftarrow FPA(\mathbf{H}, h)$ (where $\Sigma_{\mathbf{x}} = \mathbf{H}\mathbf{H}^T$ and h is the number of eigenvectors required) then the above procedure can be given as $\mathbf{w} \leftarrow FPA(\mathbf{B}, 1)$. The computational complexity for obtaining \mathbf{w} in Eq. (4) is $2dc$ and in Eq. (5) is $2dc$. Therefore, the total computational complexity is $4dc$ (in Eqs (4) and (5)).

The vector $\mathbf{w} \in \mathbb{R}^d$ is, therefore, used to transform d -dimensional space to 1-dimensional space. Let $\mathbf{x} \in$

Table 1
Gene selection procedure using fixed-point algorithm

Step 0. Define q the number of genes required and set $l = d$.

Step 1. Compute $\mathbf{w} \in \mathbb{R}^l$ using fixed-point algorithm $\mathbf{w} \leftarrow FPA(\mathbf{B}, 1)$.

Step 2. Compute z_i using Eq. (7) for $i = 1, 2, \dots, l$.

Step 3. Sort z_i in descending order; i.e., if $s = \text{sort}(z_i)$ then $s_1 > s_2 > \dots > s_l$.

Step 4. Discard least important feature corresponding to s_l . Let the cardinality of the remaining feature set be $l - 1$ and data subset be $\mathbf{X}_{l-1} \in \mathbb{R}^{l \times n}$.

Step 5. Conduct $\mathbf{X} \leftarrow \mathbf{X}_{l-1}$ and $l \leftarrow l - 1$.

Step 6. Continue Steps 1–5 until $l = q$.

\mathbf{X} be any feature vector, we have

$$y = \mathbf{w}^T \mathbf{x},$$

or

$$y = \sum_{i=1}^d w_i x_i, \quad (6)$$

where w_i and x_i are the elements of \mathbf{w} and \mathbf{x} , respectively. It can be envisaged that if $|w_i x_i| \approx 0$ (where $|\cdot|$ is the absolute value), then i th element is not contributing for the value of y in Eq. (6); i.e., it can be discarded without sacrificing much information. Therefore, we have

$$z_i = \sum_{j=1}^n |w_i x_{ij}| \quad (7)$$

where $i = 1, 2, \dots, d$. If $z_i \approx 0$, then i th feature can be discarded. Equation (7) can be applied recursively to discard unimportant features. The procedure is depicted in Table 1.

The above process will give q genes with the data subset $\mathbf{X}_q \in \mathbb{R}^{q \times n}$, which can be used by a classifier to obtain classification performance.

4. Experimentation

In this experiment we have utilized three DNA microarray gene expression datasets.¹ The description of these datasets is given as follows:

¹Most of the datasets are downloaded from the Kent Ridge Biomedical Dataset (KRBD) (<http://datam.i2r.a-star.edu.sg/datasets/krbd/>). The datasets are transformed or reformatted and made available by KRBD repository and we have used them without any further preprocessing. Some datasets which are not available on KRBD repository are downloaded and directly used from respective authors' supplement link. The URL addresses for all the datasets are given in the Reference Section.

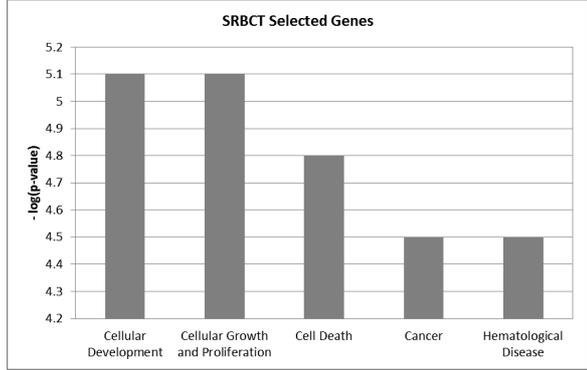


Fig. 1. Top five high level biological function on selected 150 genes of SRBCT by feature selection method based on fixed-point algorithm.

ALL dataset [7]: this dataset consists of DNA microarray gene expression data of human acute leukemia for cancer classification. Two types of acute leukemia data are provided for classification namely acute lymphoblastic leukemia (ALL) and acute myeloid leukemia (AML). The dataset is subdivided into 38 training samples and 34 test samples. The training set consists of 38 bone marrow samples (27 ALL and 11 AML) over 7129 probes. The test set consists of 34 samples with 20 ALL and 14 AML, prepared under different experimental conditions. All the samples have 7129 dimensions and all are numeric.

SRBCT dataset [10]: the small round blue-cell tumor dataset consists of 83 samples, each having 2308 genes. This is a four class classification problem. The tumors are Burkitt lymphoma (BL), the Ewing family of tumors (EWS), Neuroblastoma (NB) and Rhabdomyosarcoma (RMS). There are 63 samples for training and 20 samples for testing. The training set consists of 8, 23, 12 and 20 samples of BL, EWS, NB and RMS respectively. The test set consists of 3, 6, 6 and 5 samples of BL, EWS, NB and RMS respectively.

MLL Leukemia dataset [2]: this dataset has 3 classes namely ALL, MLL and AML. The training set contains 57 leukemia samples (20 ALL, 17 MLL and 20 AML) whereas the test set contains 15 samples (4 ALL, 3 MLL and 8 AML). The dimension of MLL dataset is 12582.

The classification performance of the proposed feature selection method has been measured on these three DNA microarray gene expression datasets. Tables 2 and 3 show classification accuracy of the proposed method compared with several other existing feature selection methods. We use J4.8 and Naïve Bayes clas-

Table 2

Classification accuracy with 150 selected genes obtained by using various feature selection methods and with J4.8 classifier on SRBCT, MLL and ALL datasets

| Feature selection methods | SRBCT (Classification accuracy) | MLL (Classification accuracy) | ALL (Classification accuracy) | Average classification accuracy (over all the 3 datasets) |
|---------------------------|------------------------------------|----------------------------------|----------------------------------|---|
| Twoing rule | 64% | 60% | 91% | 71.7 |
| Sum minority | 68% | 68% | 91% | 75.7 |
| Gini index | 64% | 60% | 91% | 71.7 |
| Sum of variances | 54% | 60% | 91% | 68.3 |
| One dimensional SVM | 54% | 60% | 91% | 68.3 |
| Fixed-point algorithm | 70% | 93% | 94% | 85.7 |

Table 3

Classification accuracy with 150 selected genes obtained by using various feature selection methods and with Naïve Bayes classifier on SRBCT, MLL and ALL datasets

| Feature selection methods | SRBCT (Classification accuracy) | MLL (Classification accuracy) | ALL (Classification accuracy) | Average classification accuracy (over all the 3 datasets) |
|---------------------------|------------------------------------|----------------------------------|----------------------------------|---|
| Twoing rule | 73% | 86% | 97% | 85.3 |
| Sum minority | 68% | 26% | 97% | 63.7 |
| Gini index | 78% | 68% | 97% | 81.0 |
| Sum of variances | 64% | 54% | 97% | 71.7 |
| One dimensional SVM | 64% | 54% | 85% | 67.7 |
| Fixed-point algorithm | 70% | 100% | 91% | 87.0 |

Table 4
Cancer functions

| Functions | <i>p</i> -value | # Selected genes |
|------------------------------|-----------------|------------------|
| Leukemia | 3.46E-05 | 13 |
| Chronic leukemia | 8.62E-05 | 8 |
| Myeloproliferative disorder | 1.51E-04 | 9 |
| Myeloid leukemia | 1.64E-04 | 8 |
| Hematologic cancer | 4.98E-04 | 14 |
| Hematological neoplasia | 5.05E-04 | 16 |
| Neuroblastoma | 1.02E-03 | 5 |
| B-cell leukemia | 1.22E-03 | 6 |
| Tumorigenesis of carcinoma | 1.32E-03 | 2 |
| Genital tumor | 1.52E-03 | 18 |
| B-cell non-Hodgkin's disease | 1.88E-03 | 6 |
| Diffuse B-cell lymphoma | 1.97E-03 | 4 |
| Prostate cancer | 2.17E-03 | 13 |
| Chronic myeloid leukemia | 2.45E-03 | 4 |
| Lymphocytic leukemia | 2.75E-03 | 7 |
| Leiomyomatosis | 2.81E-03 | 8 |
| Lymphatic node tumor | 2.99E-03 | 8 |
| Cancer | 3.86E-03 | 45 |
| Uterine leiomyoma | 3.89E-03 | 7 |
| Gliosarcoma | 4.04E-03 | 2 |

sifiers from WEKA.² The classification accuracy for SRBCT and MLL datasets is obtained from Tao et al. [21]. For Acute Leukemia dataset, the features are ranked by Rankgene program [20]. For all the datasets, we select 150 genes as done by Tao et al. [21].

It can be observed from Table 2 that the proposed method achieves highest classification accuracy (70%)

on SRBCT dataset, MLL dataset (93%) and ALL dataset (94%). The average classification accuracy of fixed-point algorithm is 85.7% which is higher than the other techniques. Furthermore, from Table 3, we can observe that average classification accuracy of fixed-point algorithm is 87% which is also higher than the other techniques. It can be concluded that the fixed-point algorithm can be applied on human cancer classification problem.

We also conducted experiments to see the biological significance of the selected features by the proposed feature selection method based on fixed-point algorithm. In order to see this, we use SRBCT data as a prototype using Ingenuity Pathway Analysis.³ The selected 150 features from the algorithm are used for this purpose. The top five high level biological functions obtained are shown in Fig. 1. In the figure, the *y*-axis denotes the negative of logarithm of *p*-values and *x*-axis denotes the high level functions. Since the cancer function is of paramount interest, we investigated them further. There are 72 cancer functions obtained from the experiment. Top 20 cancer functions with significant *p*-values are shown in Table 4. In the table, the *p*-values and the number of selected genes are depicted corresponding to the selected functions. The selected genes by the proposed method provide signifi-

²<http://www.cs.waikato.ac.nz/ml/weka/>.

³IPA, <http://www.ingenuity.com>.

cant p -values above the threshold (as specified in IPA). This shows that the features selected by the proposed method contain useful information for discriminatory purpose as well as have biological significance.

5. Conclusion

In this paper, we have presented a feature selection algorithm using fixed-point algorithm. We have shown its application in the field of human cancer classification. Three DNA microarray gene expression datasets have been utilized to see the performance of the proposed method. It was observed that the method is giving promising results. In addition, the genes selected are biologically significant as demonstrated by performing functional analysis of the genes.

References

- [1] D. Albanese, S. Merler, G. Jurman, R. Visintainer and C. Furlanello, Mlpy-high-performance Python package for predictive modeling, NIPS 08, Whistler, B.C. Canada, version 3.4.0, 2012, Software available at <http://mlpy.fbk.eu>.
- [2] S.A. Armstrong, J.E. Staunton, L.B. Silverman, R. Pieters, M.L. den Boer, M.D. Minden, S.E. Sallan, E.S. Lander, T.R. Golub and S.J. Korsmeyer, MLL translocations specify a distinct gene expression profile that distinguishes a unique leukemia, *Nature Genetics* **30** (2002), 41–47. [Data Source1: <http://sdmc.lit.org.sg/GEDatasets/Datasets.html>]; [Data Source2: http://www.broad.mit.edu/cgi-bin/cancer/publications/pub_paper.cgi?mode=view&paper_id=63].
- [3] T.-C. Chen, K. Chen, W. Liu and L.-G. Chen, Design and implementation of leading eigenvector generator for on-chip principal component analysis pike sorting system, *New Developments in Biomedical Engineering*, D. Campolo, ed., In-Tech, Jan, 2010.
- [4] T.-C. Chen, K. Chen, W. Liu and L.-G. Chen, On-chip principal component analysis with a mean pre-estimation method for spike sorting, *IEEE International Symposium on Circuits and Systems* (24–27 May 2009), 3110–3113.
- [5] T.-C. Chen, W. Liu and L.-G. Chen, VLSI architecture of leading eigenvector generation for on-chip principal component analysis spike sorting system, *IEEE EMBS 30th Annual International Conference*, Canada, (2008), 20–24.
- [6] R.O. Duda and P.E. Hart, *Pattern classification and scene analysis*, Wiley, New York, 1973.
- [7] T.R. Golub, D.K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J.P. Mesirov, H. Coller, M.L. Loh, J.R. Downing, M.A. Caligiuri, C.D. Bloomfield and E.S. Lander, Molecular classification of cancer: Class discovery and class prediction by gene expression monitoring, *Science* **286** (1999), 531–537.
- [8] M.El-B. Hazem, New fast principal component analysis for real-time face detection, *Machine Graphics & Vision International Journal* **18**(4) (2009), 405–426.
- [9] A. Hyvärinen and E. Oja, A fast fixed-point algorithm for independent component analysis, *Neural Computation* **9**(7) (1997), 1483–1492.
- [10] J. Khan, J.S. Wei, M. Ringner, L.H. Saal, M. Ladanyi, F. Westermann, F. Berthold, M. Schwab, C.R. Antonescu, C. Peterson and P.S. Meltzer, Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural network, *Nature Medicine* **7** (2001), 673–679. [Data Source: <http://research.nhgri.nih.gov/microarray/Supplement/>].
- [11] J.Z.C. Lai and T.-J. Huang, Fast global k-means clustering using cluster member and inequality, *Pattern Recognition* **43** (2010), 1954–1963.
- [12] K.K. Paliwal and A. Sharma, Improved direct LDA and its application to DNA microarray gene expression data, *Pattern Recognition Letters* **31**(16) (2010), 2489–2492.
- [13] R.C. Qiu, Z. Chen, N. Guo, Y. Song, P. Zhang, H. Li and L. Lai, Towards a real-time cognitive radio network testbed: Architecture, hardware platform, and application to smart grid, *Fifth IEEE Workshop on Networking Technologies for Software Defined Radio (SDR) Networks*, Boston, MA, USA (21–21 June 2010).
- [14] K. Ramesha and K.B. Raja, Face recognition system using discrete wavelet transform and fast PCA, *Information Technology and Mobile Communication, Communications in Computer and Information Science* **147**(1) (2011), 13–18.
- [15] I. Sajid, M.M. Ahmed and I. Taj, Design and implementation of a face recognition system using fast PCA, *International Symposium on Computer Science and its Applications, CSA'08*, (2008), 126–130.
- [16] A. Sharma and K.K. Paliwal, Fast principal component analysis using fixed-point algorithm, *Pattern Recognition Letters* **28**(10) (2007), 1151–1155.
- [17] A. Sharma and K.K. Paliwal, A gradient linear discriminant analysis for small sample sized problem, *Neural Processing Letters* **27**(1) (2008), 17–24.
- [18] A. Sharma and K.K. Paliwal, A new perspective to null linear discriminant analysis method and its fast implementation using random matrix multiplication with scatter matrices, *Pattern Recognition* **45** (2012), 2205–2213.
- [19] W. Shi and Y.-F. Guo, Nonlinear component analysis for large-scale data set using fixed-point algorithm, *Advances in Neural Networks, Lecture Notes in Computer Science* **5553** (2009), 144–151.
- [20] Y. Su, T.M. Murali, V. Pavlovic and S. Kasif, RankGene: Identification of diagnostic genes based on expression data, *Bioinformatics* (2003), 1578–1579.
- [21] L. Tao, C. Zhang and M. Ogihara, A comparative study of feature selection and multiclass classification methods for tissue classification based on gene expression, *Bioinformatics* **20**(14) (2004), 2429–2437.
- [22] J. Wang, A. Barreto, N. Rishe and J. Andrian, A fast incremental multilinear principal component analysis algorithm, *International Journal of Innovative Computing, Information and Control* **7**(10) (2011), 6019–6040.
- [23] C. Yang, L. Wang and J. Feng, A novel margin based algorithm for feature extraction, *New Generation Computing* **27** (2009), 285–305.
- [24] P. Zhang, R. Qiu and N. Guo, Demonstration of spectrum sensing with blindly learned features, *IEEE Communication Letters* **15**(5) (2011), 548–550.